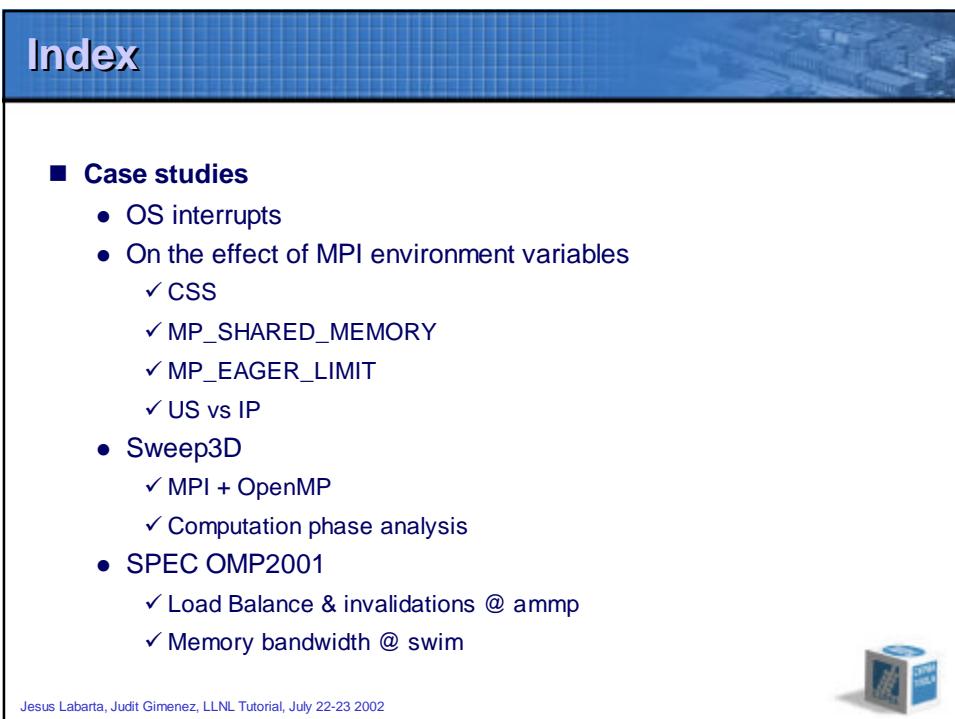




The slide features a background image of a city skyline at night. In the top left corner, there is a small blue icon of a computer monitor with three parallel bars on its screen. The title "Paraver Advanced" is centered in large, dark blue serif font. Below the title, the names "Jesús Labarta, Judit Giménez", "Jordi Caubet, Francesc Escale", and the acronym "CEPBA-UPC" are listed in a smaller, dark blue sans-serif font. At the bottom of the slide, there is a horizontal menu bar with four items: "Technology Transfer", "Research", "Training", and "Mobility of Researchers". Underneath each of these main items are two sub-items: "User Support", "Education", "HPC Facilities", and "Parallel Expertise".



The page has a blue header bar with the word "Index" in white. The main content area contains a section titled "Case studies" with a black square bullet point. This section lists several case studies with associated checkmarks:

- **Case studies**
 - OS interrupts
 - On the effect of MPI environment variables
 - ✓ CSS
 - ✓ MP_SHARED_MEMORY
 - ✓ MP_EAGER_LIMIT
 - ✓ US vs IP
 - Sweep3D
 - ✓ MPI + OpenMP
 - ✓ Computation phase analysis
 - SPEC OMP2001
 - ✓ Load Balance & invalidations @ ammp
 - ✓ Memory bandwidth @ swim

At the bottom left, there is a small note: "Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002". At the bottom right, there is a small blue cube icon with the text "CEPBA-UPC" on it.

OS internals

Jesús Labarta, Judit Giménez, Jordi Caubet
CEPBA-UPC

Technology Transfer | Research | Training | Mobility of Researchers
User Support | Education | HPC Facilities | Parallel Expertise

Tracing overhead

```
DO I=1,10000
  call ompitrace_eventandcounters(1000,I)
  call ompitrace_eventandcounters(1000,0)
END DO
```

- OpenMP
 - 0 – 100 µs

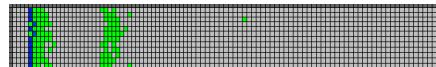
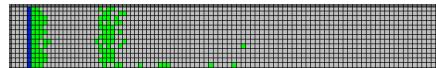
Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002

Tracing overhead

```
DO I=1,10000
call ompitrace_eventandcounters(1000,I)
call ompitrace_eventandcounters(1000,0)
END DO
```

■ MPI

- 0 – 100 µs



Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002

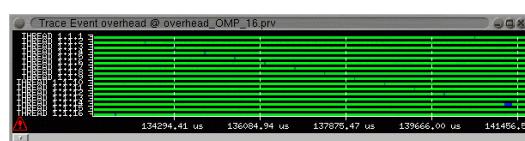
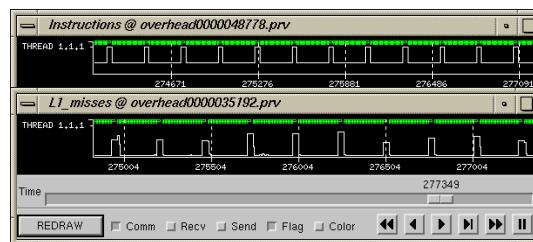
System interrupts

■ Observable

- Per process kernel lock in access to hardware counters information

- Quantum
- Cost

- Skew between CPUs



Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002

On the effect of environment variables

Jesús Labarta, Judit Giménez, Jordi Caubet
CEPBA-UPC

Technology Transfer | Research | Training | Mobility of Researchers |
User Support | Education | HPC Facilities | Parallel Expertise

Effect of environment variables

- 16 tasks
- Environment
 - MP_SHARED_MEMORY
 - MP_EAGER_LIMIT
 - ✓ 4k / 64k
 - MP_CSS_INTERRUPT
 - IP / US
 - Number of nodes
 - ✓ 1, 2, 4, 8

LU

BT

IS

Sweep3D

Jesus Labarta, Judit Giménez, LLNL Tutorial, July 22-23 2002

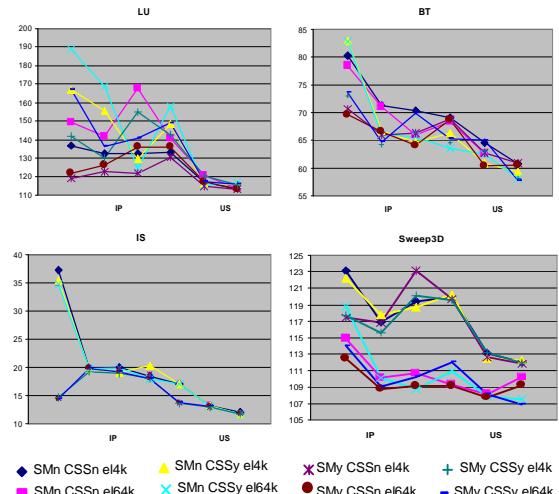
Effect of environment variables

- Often <10%, sometimes significant

- Good combination:

- unpredictable
- Autonomic Computing needed here too

- US reduces influence of others

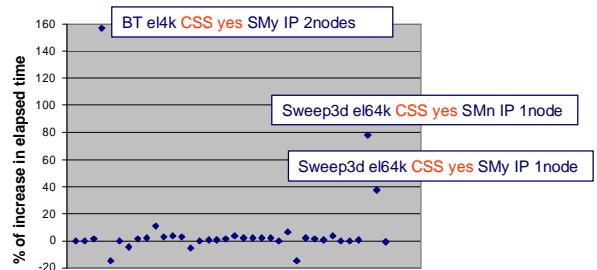


Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002

Tracing overhead

- Comparison of elapsed time when tracing vs. un-instrumented

- Generally small difference (<5%)
- CSS responsible of outliers



Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002

Studies

- Effect of MP_SHARED_MEMORY
- Effect of MP_EAGER_LIMIT
- Comparison between US and IP
- MP_CSS_INTERRUPT
- 2 nodes better than 1 node
- US vs. SM
- ...
- ...
-

8

Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002

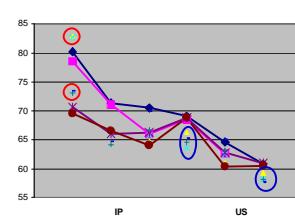


CSS

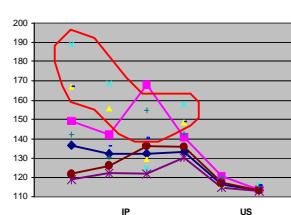
- Usually bad

- May be good

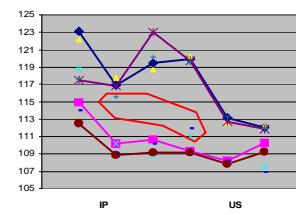
BT



LU



Sweep3D



Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002



CSS:

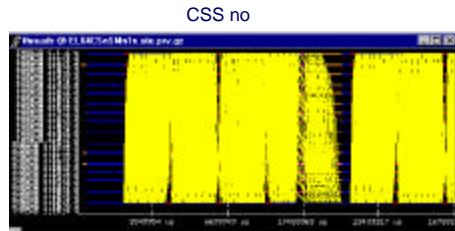
... sweep3d

■ The bad

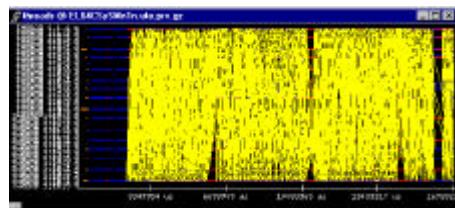
■ OS scheduling

- Conflict between application and MPI implementation threads

■ Ute traces ? Paraver



Sweep3d, SMy, 1 node



CSS yes

Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002

CSS:

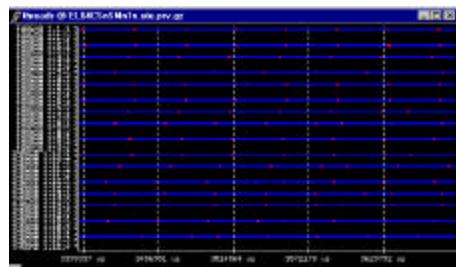
... sweep3d

■ The bad

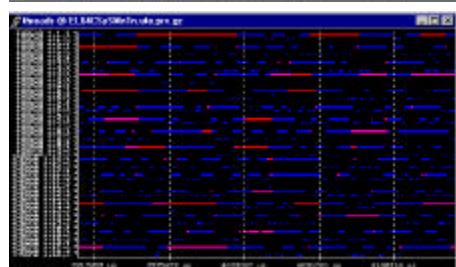
■ OS scheduling

- Is yielding used?

CSS no



CSS yes



Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002

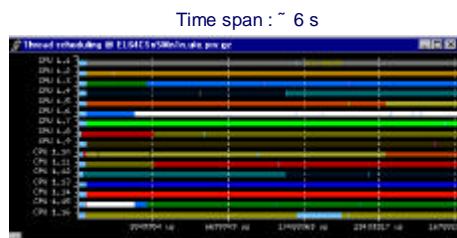
CSS:

... sweep3d

■ The bad

- Process migrations
- Can internal threads be bound to their corresponding master thread?

CSS no



CSS yes



Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002



CSS:

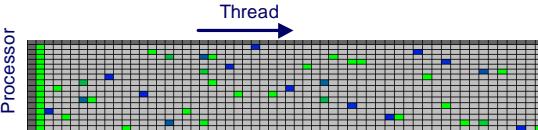
... sweep3d

■ Time spent by each thread on each processor

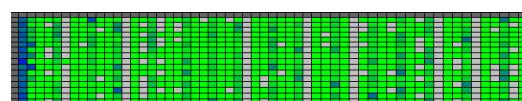
- Visible: Thread to processor mapping if CSS ==no
- Curious: Single thread taking more % of CPUtime at each processor if CSS == yes: Idle processor (15%)

Processor

Thread →



CSS no



CSS yes

Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002

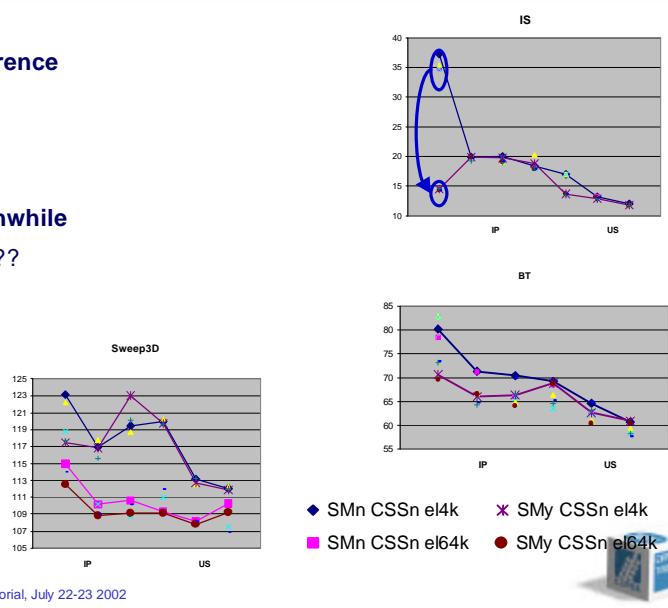


MP_SHARED_MEMORY

- Important difference

- Or not

- Generally worthwhile
 - Default NO ??



Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002

MP_SHARED_MEMORY:

... IS

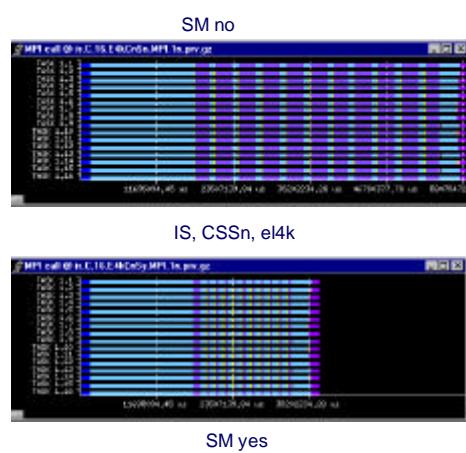
- Load MPI_call configuration

- Important difference

- Need to quantify

- Shared memory: first time slower

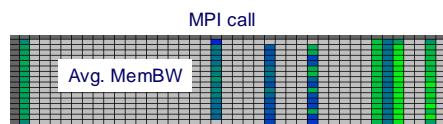
- Significant influence in average



Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002

MP_SHARED_MEMORY: ... IS

- 2D analysis configuration file
- A given environment set up may be good for one call and not so good for other
- Statistics
 - Care with selected area
 - May be misleading
 - ✓ Low miss ratio not always synonym of good



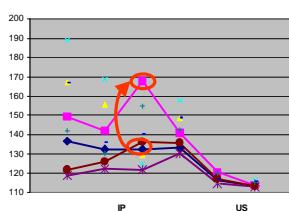
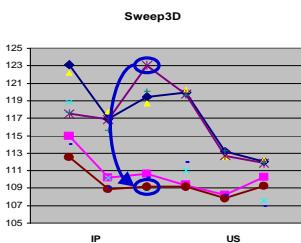
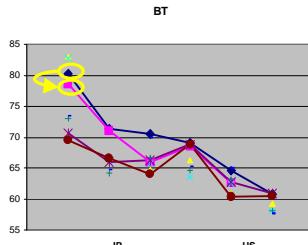
	SM yes	SM no	Ratio
time (s)	14.46	37.25	0.39
MPI_Altally duration (us)	441314	2488005	0.18
MPI_Altoall duration (us)	520	8334	0.06
MPI_Allreduce duration (us)	205017	186966	1.10

Avg. Values @ MPI_Altally	SM yes	SM no	Ratio
Stores	9421214	29634927	0.32
Loads	11632523	66742482	0.17
Mem Ops mix (%)	74	43.6	1.70
L2 miss ratio (%)	4.1	1.6	2.56
Mem BW (MB/s)	375.8	94.2	3.99
CPU BW (MB/s)	440	310	1.42
IPC	0.2	0.24	0.83

Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002

MP_EAGER_LIMIT

- Increasing eager limit we expect
 - Not much difference
 - An improvement
- but may get bad !!!



◆ SMn CSSn el4k ✕ SMy CSSn el4k
 ■ SMn CSSn el64k ● SMy CSSn el64k

Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002

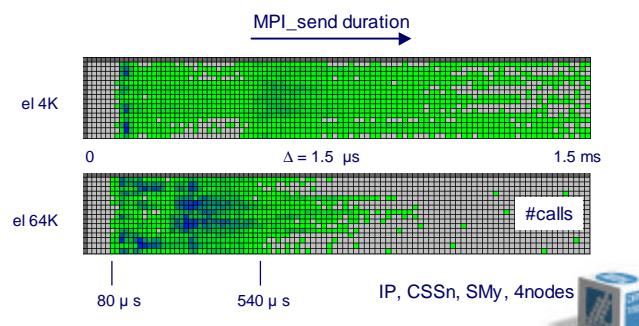
MP_EAGER_LIMIT: ... sweep3d

■ Expectation:

- Sufficient Eager limit (64K) → sends do not wait

■ Time distribution of MPI_send

- Fast calls approximately equal
- Slow calls faster

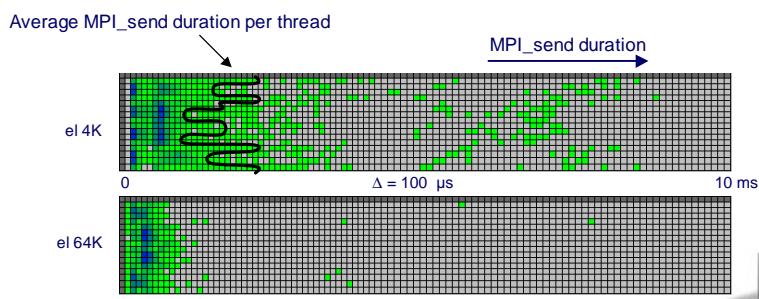


Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002

MP_EAGER_LIMIT: ... sweep3d

■ A view from the heights

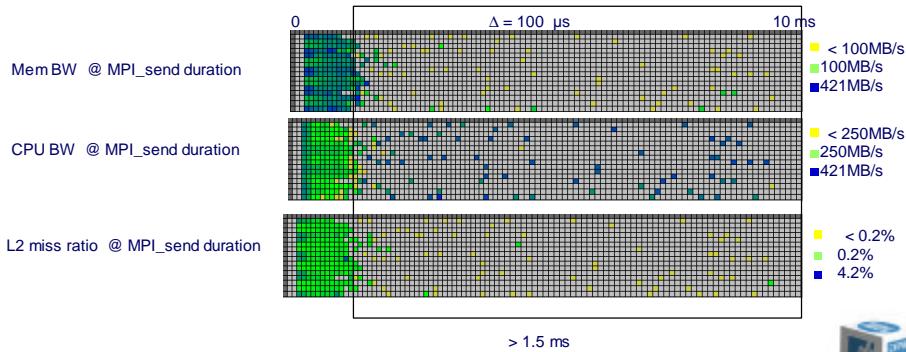
- Some very slow sends
- Not uniform distribution
 - ✓ Time distribution
 - ✓ Per thread distribution



Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002

MP_EAGER_LIMIT: ... sweep3d

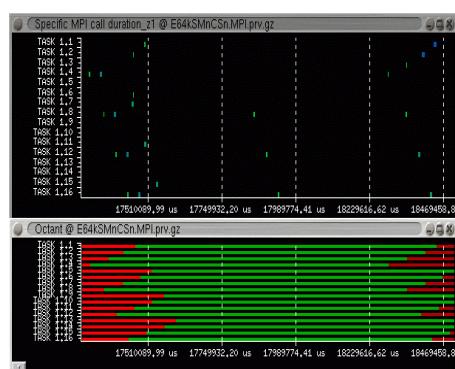
- Some very long sends even with large eager limit
 - Configuration 1node



Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002

MP_EAGER_LIMIT: ... sweep3d

- Identification
 - Few cases
 - ✓ difficult based on direct navigation state timeline
 - Specific view
 - ✓ MPI sends of duration >1.5ms



■ Correlation

- With source
- With iteration

Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002

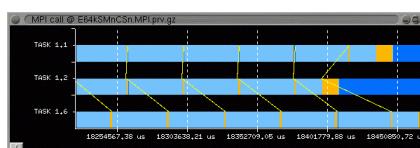
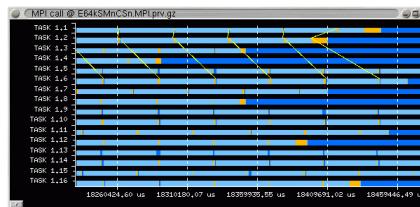
MP_EAGER_LIMIT: ... sweep3d

■ Focus View

- Filter Displayed communications
- Hide non relevant threads

■ Interpretation

- MP_EAGER_LIMIT ???
 - ✓ Several consecutive sends
 - ✓ What if previous one not yet received ?



Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002

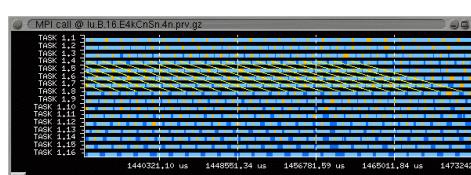
MP_EAGER_LIMIT: ... LU

■ Increasing eager limit is bad !!!

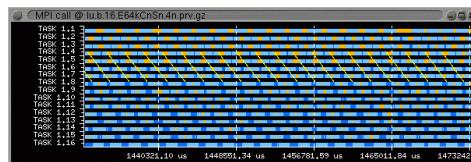
- For some calls
- Even if all messages <4K
- Flow control issue

	4K	64K	Ratio
time (s)	136.20	167.90	0.81
MPI_send duration (us)	162.4	255.8	0.63
MPI_recv duration (us)	225.6	273.2	0.83
MPI_irecv duration (us)	54.4	54.6	1.00
MPI_Wait duration (us)	7260	4544	1.60

4K



64K

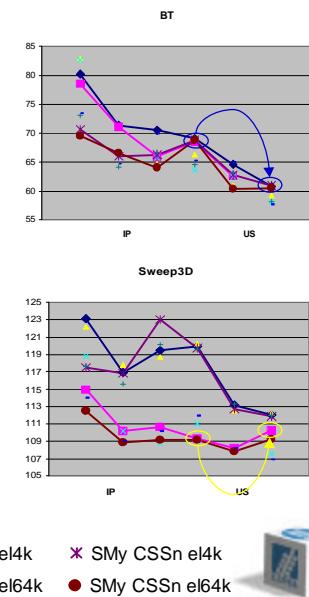


CSS no, SM no, 4 nodes

Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002

US vs IP

- Improves for BT and IS
- Very little for LU, nothing for Sweep3d if proper eager limit

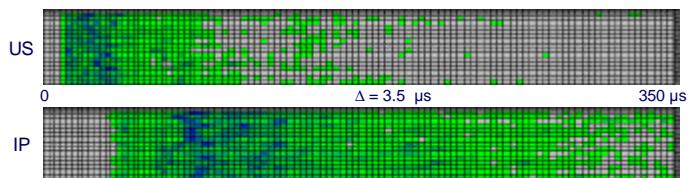


Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002

US vs IP:

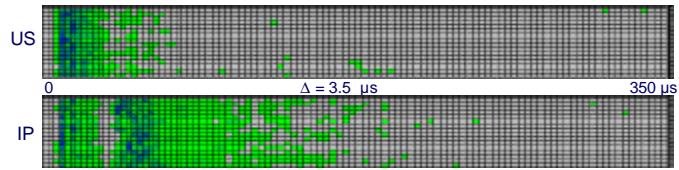
... BT

- **I send**
 - Speed
 - US less variance



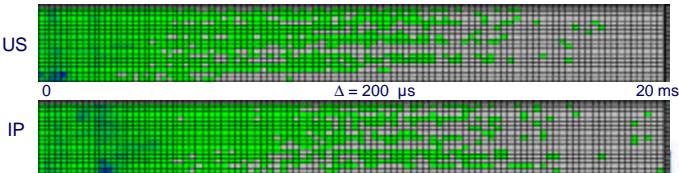
- **I recv**

- Uni/Bi modal



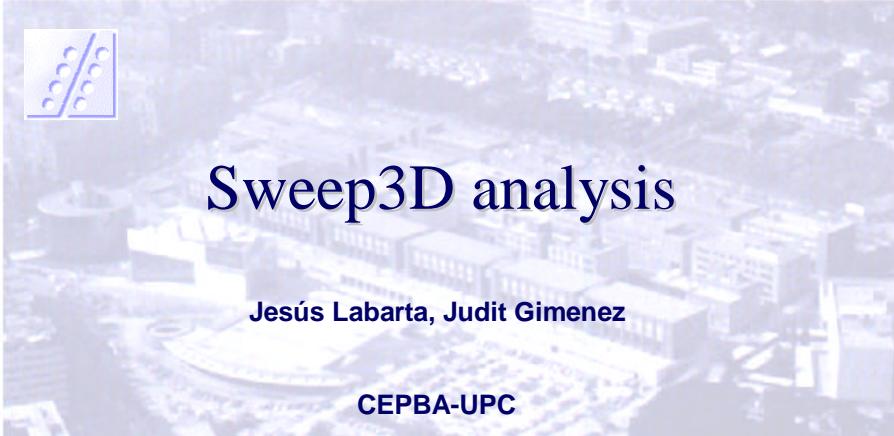
- **Wait**

- Variance



Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002

BT, CSSn, SMn, e4k, 8 nodes



Sweep3D analysis

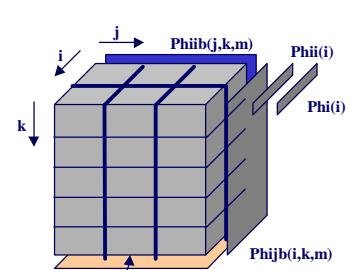
Jesús Labarta, Judit Giménez

CEPBA-UPC

Technology Transfer Research Training Mobility of Researchers
User Support Education HPC Facilities Parallel Expertise

MPI Parallelization

- Grid Partitioning of i, j plane
- One column per MPI task
- Reduction on each direction:
 - Communicate on i and j
 - Pipeline k dimension
 - + overlap consecutive sweeps if possible
- Input file parameters
 - Blocks on i and j directions
 - Pipelining block



flux(i,j,k,n), face(i,j,k,n), src(i,j,k,n)

MPI data parallelization
Flux, face,src: DISTRIBUTED
phi_{ii}, phi: PRIVATIZED
phi_{ikb}: DISTRIBUTED
phi_{ijkb},phi_{iiib}: DISTRIBUTED&REPLICATED
=> Communication

Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002

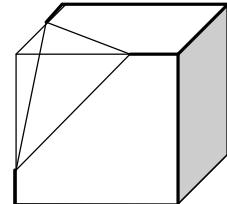
OpenMP Parallelization

■ Thought for mixed mode (MPI + OpenMP)

- Just core computational loop

■ Reduction on each direction

- Compute along diagonal wavefronts



```
DO idiag
  DO jkm =1, #points in wavefront
    j,k,m=f(idiag,jkm)
    DO n,i      ! phi, src
    DO i        ! phijb, phikb, phii, phi
    DO n,i      ! flux, phi
    DO i        ! face, phii, phijb, phikb
```

Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002



MPI vs. OpenMP: some numbers

■ Problem

- size: 50^3
- k plane pipelining: 10

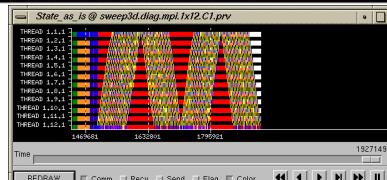
NB Domains	OpenMP time	Decomposition	MPI time
6	7.78	1x6	3.97
		2x3	3.61
		3x2	3.71
		6x1	4.47
12	6.55	1x12	3.50
		2x6	2.74
		3x4	2.21
		4x3	2.25
		6x2	2.97
		12x1	3.98

Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002

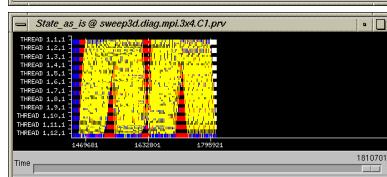


MPI : decomposition effect

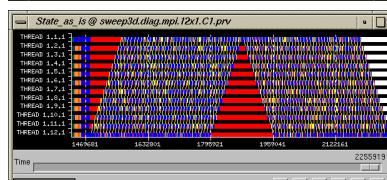
■ 1 x 12



■ 3 x 4



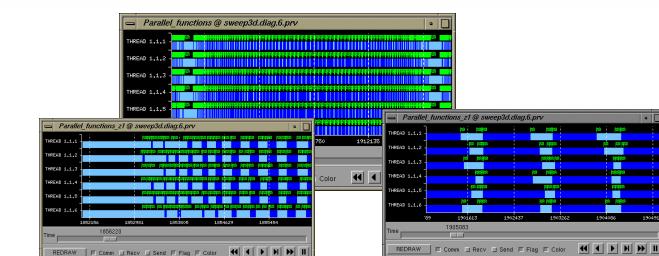
■ 12 x 1



Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002



Pure OpenMP



```
DO idiag
:
DO jkm = 1, #points in wavefront
j,k,m=f(idiag,jkm)
DO n,i      ! phi, src
DO i       ! phijb, phikb, phii, phi
DO n,i      ! flux, phi
DO i       ! face, phii, phijb, phikb
```

Parallel

Computation:
Complex
Overhead
Triangular trip count
OpenMP RTL overhead
Invalidation traffic

Computation:
Complex
Overhead

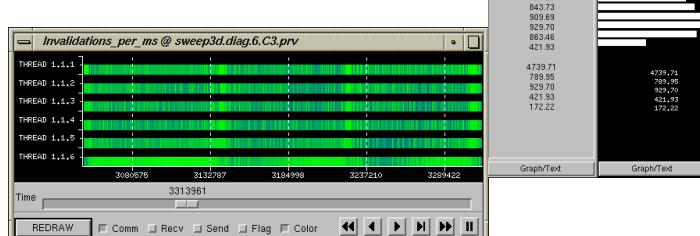
Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002



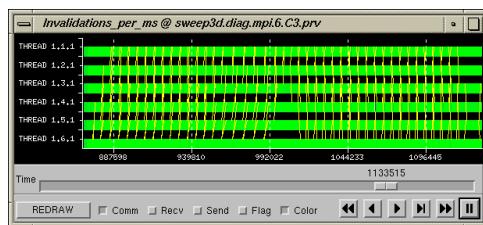
Pure OpenMP

■ Invalidations

- OpenMP



- MPI



Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002



Alternative structures

```
DO idiag
  DO jkm =1, #points in wavefront
    j,k,m=f(idiag,jkm)
    DO n,i      ! phi, src
    DO i        ! phijb, phikb, phii, phi
    DO n,i      ! flux, phi
    DO i        ! face, phii, phijb, phikb

    DO m
      DO k
        DO j
          DO n,i      ! phi, src
          DO i        ! phijb, phikb, phii, phi
          DO n,i      ! flux, phi
          DO i        ! face, phii, phijb, phikb
```

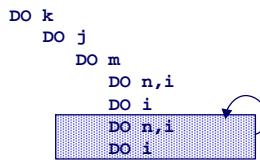
mkji form
in the distributed source

Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002

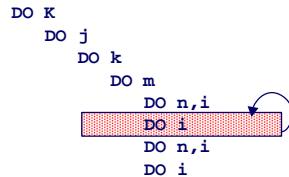


Alternative forms

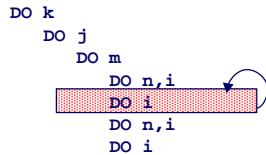
■ mkji



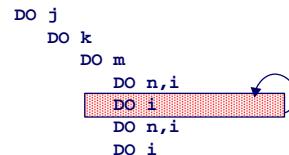
■ Kjkmi



■ kjmi



■ jkmi



■ . . .

Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002



OpenMP: some numbers

■ Problem

- size: 50^3
- k plane pipelining: 10

Elapsed Time	OMP_NUM_THREADS																
	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
ccrit	24.40	28.26	24.41	26.84	26.47	29.28	30.34										30.43
ccpipe	24.60	25.63	18.45	13.01	12.53	10.06	7.67										7.76
diag	16.21	17.28	13.09	11.40	9.64	8.50	7.78										6.55
jkmi	12.95	14.86	10.01	7.35	5.82	4.89	4.34	3.80	3.62	3.38	3.09	3.04	2.88	2.69	2.64	2.53	
Kjkmi	12.94	14.91	8.47	6.35	4.91	4.24	3.58	3.46	2.90	2.81	2.78	2.65	2.29	2.22	2.16	2.19	2.15

Invalidations

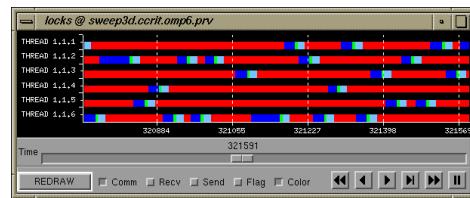
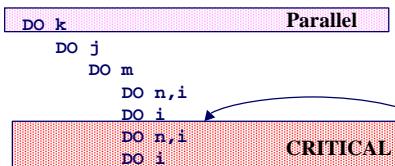
Instruction count overhead?

Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002



OpenMP : contention on locks

■ Version ccrit, 6 Threads



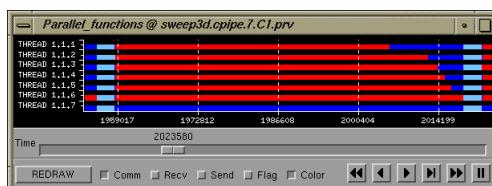
Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002



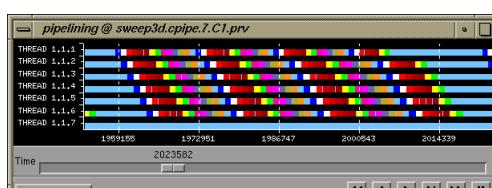
OpenMP: Insufficient parallelism

■ Version cpipe, 7 threads

- outer iteration count: 6
- parallel + worksharing



- Internal pipelined iteration



Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002



Comparing modes

Single mode better

MPI	partition	Seq	OpenMP				
			1	2	4	8	16
0	-				66,33	37,39	22,77
1	1x1				66,99	38,21	23,55
2	1x2			79,34	50,90	36,58	
	2x1			89,60	57,52	42,36	
4	1x4	66,20	71,67	41,92	28,24		
	2x2	79,72	86,71	52,01	34,60		
	4x1	76,04	82,04	54,51	37,05		
8	1x8	37,83	40,27	24,90			
	2x4	42,95	46,01	28,84			
	4x2	46,39	50,74	33,79			
	8x1	44,71	50,91	36,34			
16	1x16	23,12	24,82				
	4x4	24,75	26,86				
	16x1	28,40	32,60				

Decomposition effect

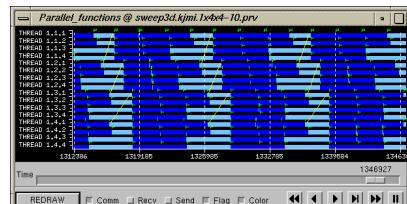
Scheduling Interference

Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002

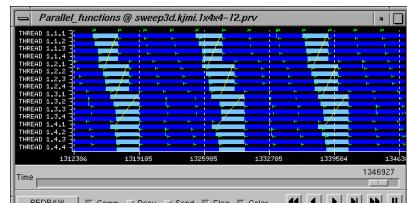
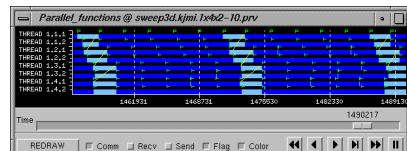


Mixed mode: scheduling interference

- MPI (4 tasks) + kjmi
 - 4 threads, k pipeline=10



- 2 threads, k pipeline=10
- 4 threads, k pipeline=12
 - ✓ Less K iterations



Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002

Computation distribution histograms

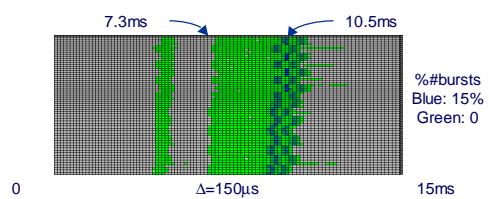
- Sweep3d (300 x 300 x 125) , 8x8x1 (MPI_I x MPI_J x OpenMP)

- 2D

- Cw: core loop duration
- Statistic: % #bursts

- Comments

- Load imbalance ?



Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002

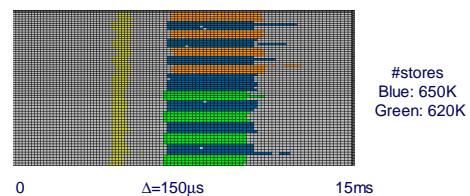


Computation distribution histograms

- 8x8x1

- 2D

- Cw: core loop duration
- Statistic: average value
- Dw:#stores



- Comments

- Load imbalance: stores = f (CPU)
✓ $300 = 37 * 8 + 4$
- Precision of measurements (consistent detection of small percentual variations)
- #Stores ≠ f(time) ?

Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002

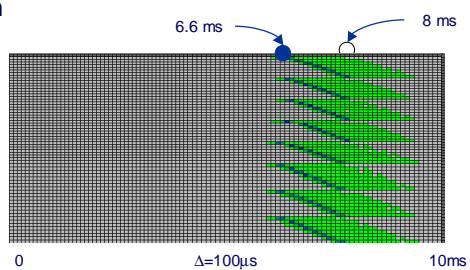


Computation distribution histograms

- Sweep3d (300 x 300 x 125) ,1x8x8 (MPI_I x MPI_J x OpenMP)

- 2D

- Cw: core loop duration
- Statistic: #bursts



- Comments

- Pipeline
- ✓ 175 μs

Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002

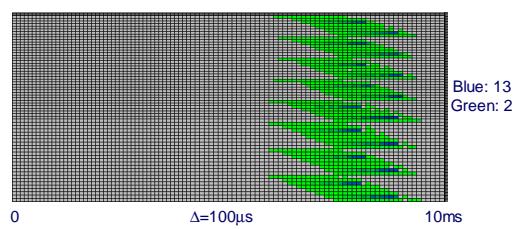


Computation distribution histograms

- 1x8x8

- 2D

- Cw: core loop duration
- Statistic: average value
- Dw: CPU/memory BW ratio



- Comments

- Generally low reuse
- There is structure
- ✓ Interpretation ?

Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002



Performance indices correlation

■ Example

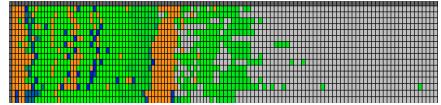
- L2 misses vs Memory BW



- IPC vs Memory BW



Loads vs Memory BW



Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002



Analysis of the SPEC OMP2001 Benchmarks with Paraver

Jesús Labarta, Judit Gimenez

CEPBA-UPC

Technology Transfer

User Support

Research

Education

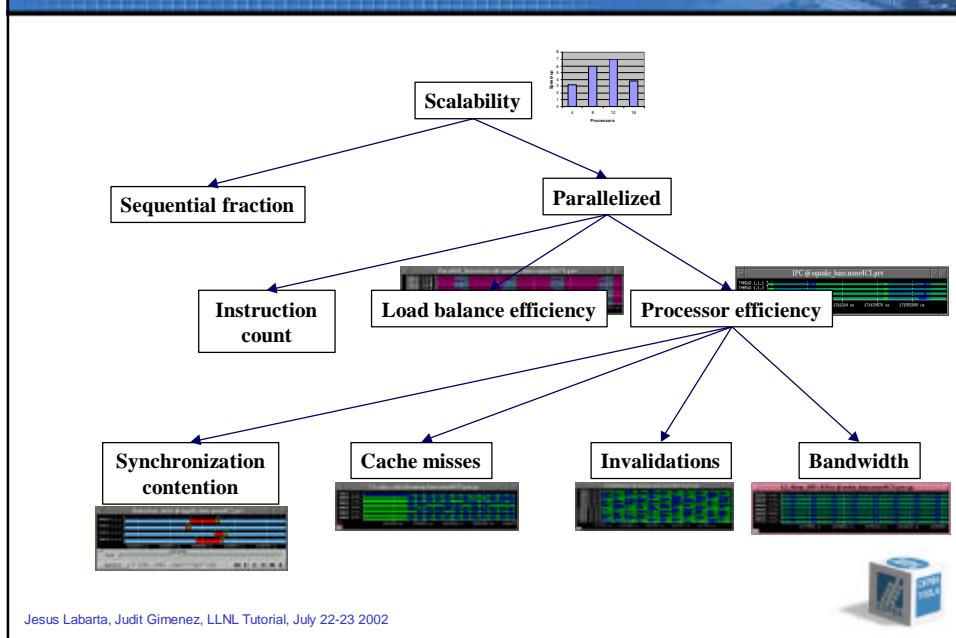
Training

HPC Facilities

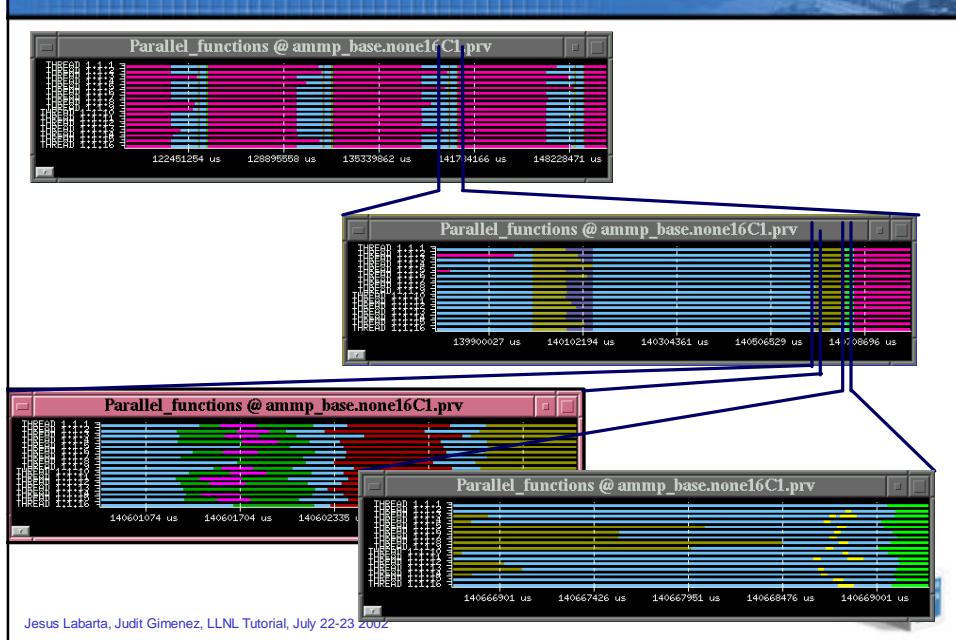
Mobility of Researchers

Parallel Expertise

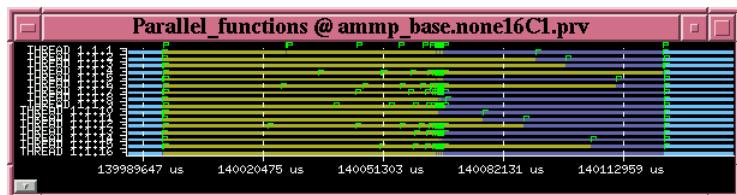
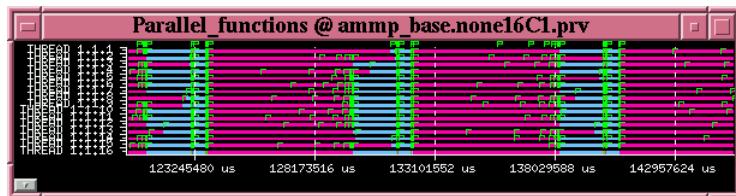
Analysis methodology



Ammp: load balance



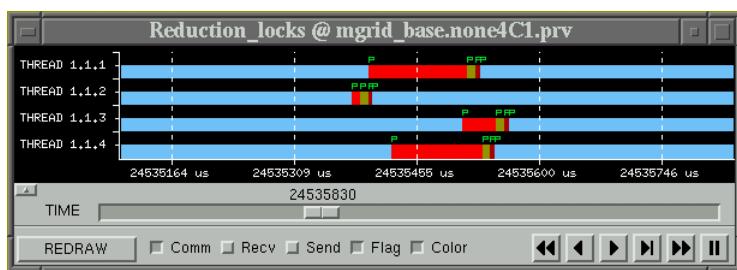
Ammp: loop scheduling



Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002



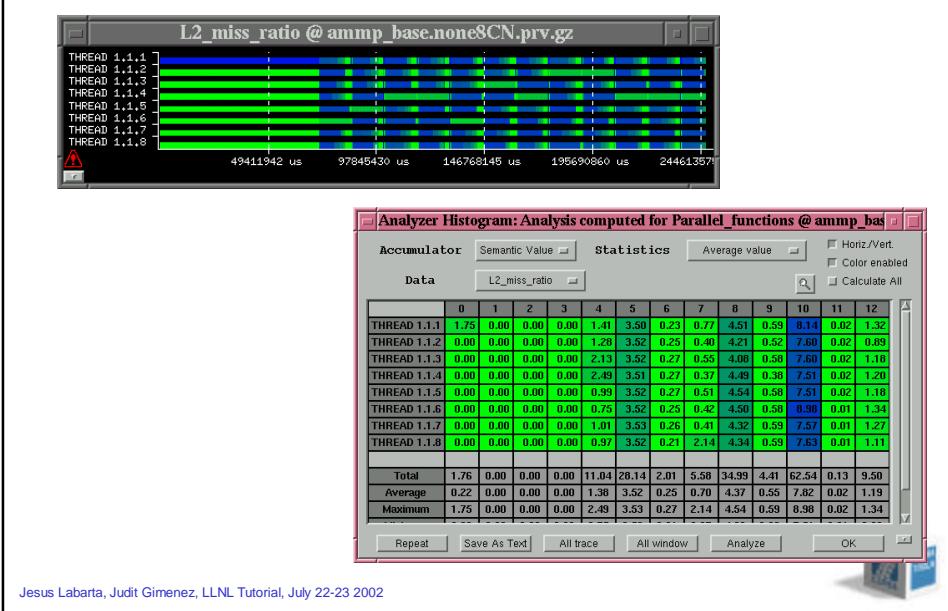
Synchronization



Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002

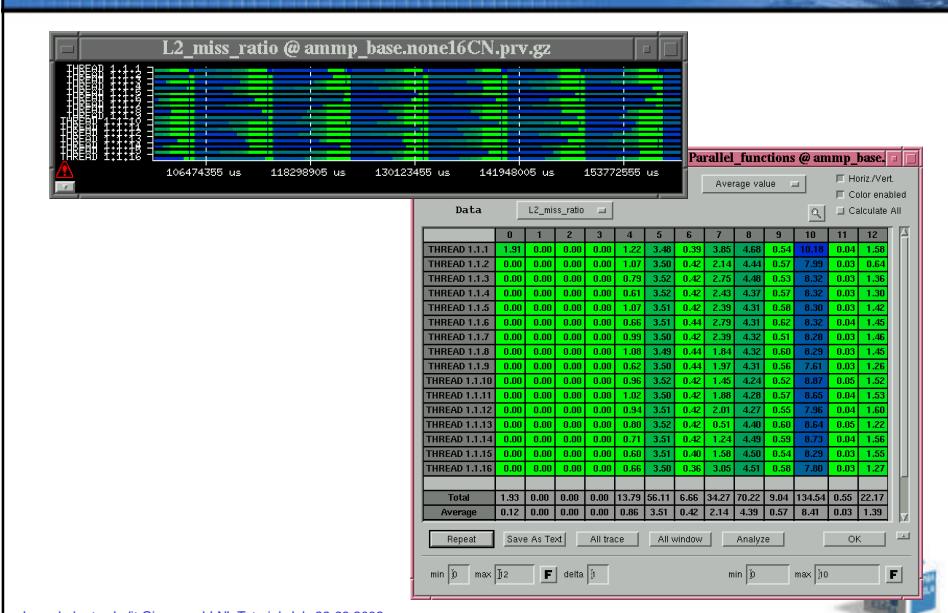


Ammp: L2 miss ratio



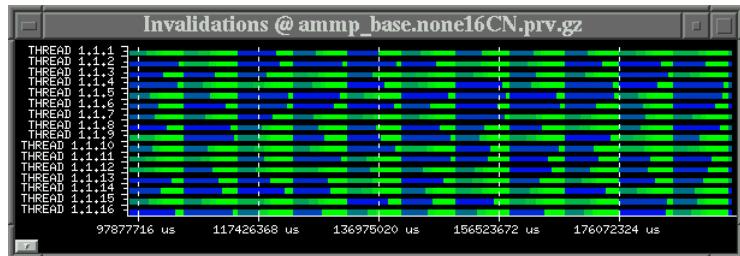
Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002

Ammp: L2 miss ratio



Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002

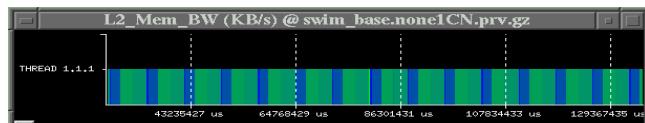
Invalidations



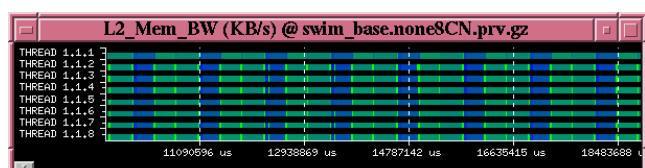
Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002



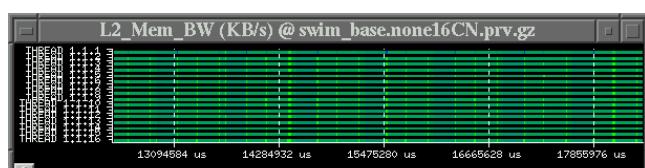
Bandwidth



Larger Y scale

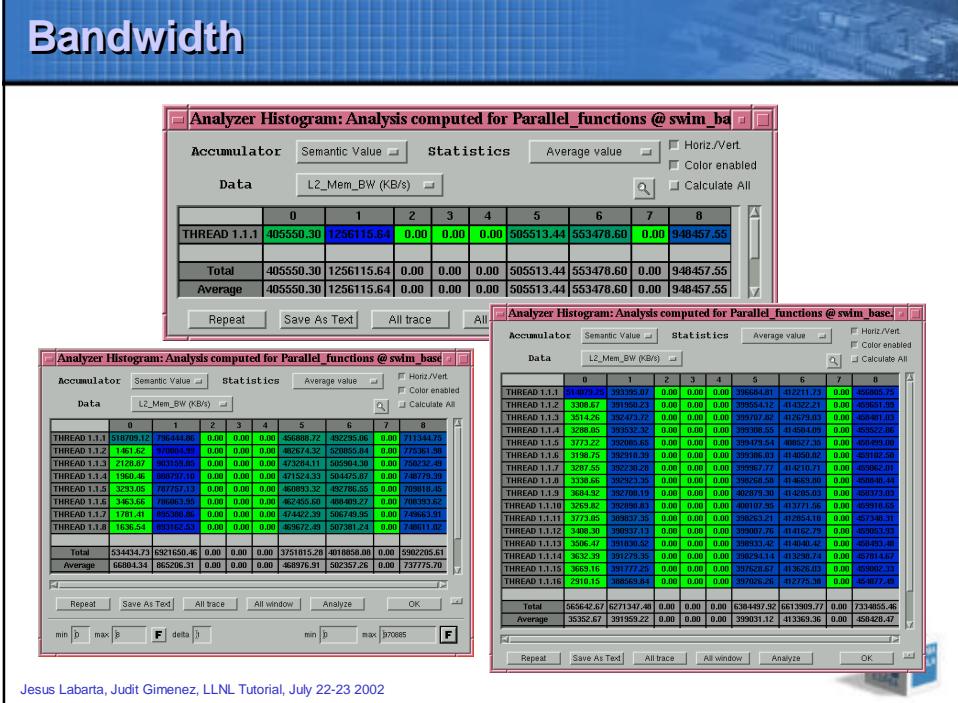


Same Y scale



Jesus Labarta, Judit Gimenez, LLNL Tutorial, July 22-23 2002





Jesús Labarta, Judit Giménez, LLNL Tutorial, July 22-23 2002